

Example 6.3: Assess a student's performance during his course of study and predict whether a student will get a job offer or not in his final year of the course. The training dataset T consists of 10 data instances with attributes such as 'CGPA', 'Interactiveness', 'Practical Knowledge' and 'Communication Skills' as shown in Table 6.3. The target class attribute is the 'Job Offer'.

Table 6.3: Training Dataset T

S.No.	CGPA	Interactiveness	Practical Knowledge	Communication Skills	Job Offer
1.	≥ 9	Yes	Very good	Good	Yes
2.	≥ 8	No	Good	Moderate	Yes
3.	≥ 9	No	Average	Poor	No
4.	< 8	No	Average	Good	No
5.	≥ 8	Yes	Good	Moderate	Yes
6.	≥ 9	Yes	Good	Moderate	Yes
7.	< 8	Yes	Good	Poor	No
8.	≥ 9	No	Very good	Good	Yes
9.	≥ 8	Yes	Good	Good	Yes
10.	≥ 8	Yes	Average	Good	Yes

Solution:

Step 1:

Calculate the Entropy for the target class 'Job Offer'.

$$\text{Entropy_Info}(\text{Target Attribute} = \text{Job Offer}) = \text{Entropy_Info}(7, 3) =$$

$$= -\left[\frac{7}{10} \log_2 \frac{7}{10} + \frac{3}{10} \log_2 \frac{3}{10} \right] = -(-0.3599 + -0.5208) = 0.8807$$

Iteration 1:

Step 2:

Calculate the Entropy_Info and Gain(Information_Gain) for each of the attribute in the training dataset.

Table 6.4 shows the number of data instances classified with Job Offer as Yes or No for the attribute CGPA.

Table 6.4: Entropy Information for CGPA

CGPA	Job Offer = Yes	Job Offer = No	Total	Entropy
≥ 9	3	1	4	
≥ 8	4	0	4	0
< 8	0	2	2	0

$$\text{Entropy_Info}(T, \text{CGPA})$$

$$= \frac{4}{10} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] + \frac{4}{10} \left[-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right] + \frac{2}{10} \left[-\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} \right]$$

$$= \frac{4}{10} (0.3111 + 0.4997) + 0 + 0$$

$$= 0.3243$$

$$\text{Gain (CGPA)} = 0.8807 - 0.3243$$

$$= 0.5564$$

Table 6.5 shows the number of data instances classified with Job Offer as Yes or No for attribute Interactiveness.

Table 6.5: Entropy Information for Interactiveness

Interactiveness	Job Offer = Yes	Job Offer = No	Total	Entropy
YES	5	1	6	
NO	2	2	4	

$$\text{Entropy_Info}(T, \text{Interactiveness}) = \frac{6}{10} \left[-\frac{5}{6} \log_2 \frac{5}{6} - \frac{1}{6} \log_2 \frac{1}{6} \right] + \frac{4}{10} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right]$$

$$= \frac{6}{10} (0.2191 + 0.4306) + \frac{4}{10} (0.4997 + 0.4997)$$

$$= 0.3898 + 0.3998 = 0.7896$$

$$\text{Gain(Interactiveness)} = 0.8807 - 0.7896$$

$$= 0.0911$$

Table 6.6 shows the number of data instances classified with Job Offer as Yes or No for attribute Practical Knowledge.

Table 6.6: Entropy Information for Practical Knowledge

Practical Knowledge	Job Offer = Yes	Job Offer = No	Total	Entropy
Very Good	2	0	2	0
Average	1	2	3	
Good	4	1	5	

Entropy_Info(T, Practical Knowledge)

$$\begin{aligned}
 &= \frac{2}{10} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right] + \frac{3}{10} \left[-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right] + \frac{5}{10} \left[-\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right] \\
 &= \frac{2}{10}(0) + \frac{3}{10}(0.5280 + 0.3897) + \frac{5}{10}(0.2574 + 0.4641) \\
 &= 0 + 0.2753 + 0.3608 \\
 &= 0.6361
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain(Practical Knowledge)} &= 0.8807 - 0.6361 \\
 &= 0.2446
 \end{aligned}$$

Table 6.7 shows the number of data instances classified with Job Offer as Yes or No for the attribute Communication Skills.

Table 6.7: Entropy Information for Communication Skills

Communication Skills	Job Offer = Yes	Job Offer = No	Total
Good	4	1	5
Moderate	3	0	3
Poor	0	2	2

Entropy_Info(T, Communication Skills)

$$\begin{aligned}
 &= \frac{5}{10} \left[-\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right] + \frac{3}{10} \left[-\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} \right] + \frac{2}{10} \left[-\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} \right] \\
 &= \frac{5}{10}(0.5280 + 0.3897) + \frac{3}{10}(0) + \frac{2}{10}(0) \\
 &= 0.3609
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain(Communication Skills)} &= 0.8813 - 0.36096 \\
 &= 0.5203
 \end{aligned}$$

The Gain calculated for all the attributes is shown in Table 6.8:

Table 6.8: Gain

Attributes	Gain
CGPA	0.5564
Interactiveness	0.0911
Practical Knowledge	0.2246
Communication Skills	0.5203

Step 3: From Table 6.8, choose the attribute for which entropy is minimum and therefore the gain is maximum as the best split attribute.

The best split attribute is CGPA since it has the maximum gain. So, we choose CGPA as the root node. There are three distinct values for CGPA with outcomes ≥ 9 , ≥ 8 and < 8 . The entropy value is 0 for ≥ 8 and < 8 with all instances classified as Job Offer = Yes for ≥ 8 and Job Offer = No for < 8 . Hence, both ≥ 8 and < 8 end up in a leaf node. The tree grows with the subset of instances with CGPA ≥ 9 as shown in Figure 6.3.

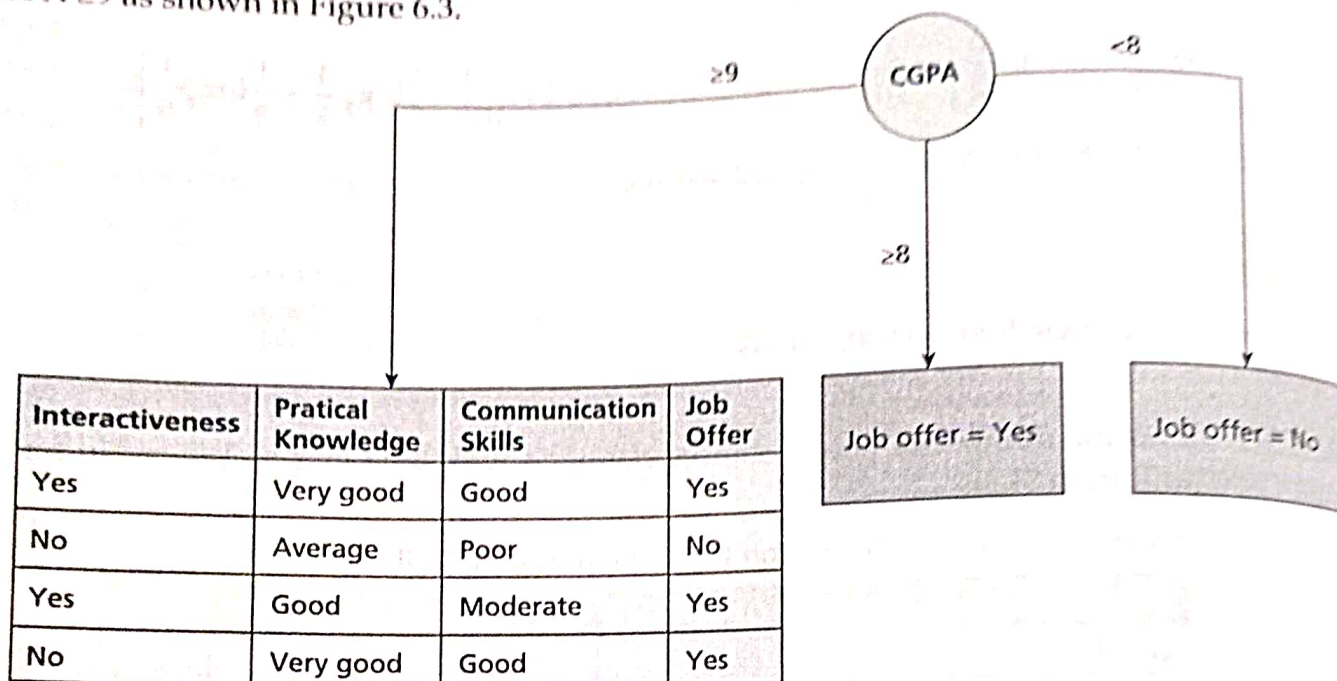


Figure 6.3: Decision Tree After Iteration 1

Now, continue the same process for the subset of data instances branched with CGPA ≥ 9 .

Iteration 2:

In this iteration, the same process of computing the Entropy_Info and Gain are repeated with the subset of training set. The subset consists of 4 data instances as shown in the above Figure 6.3.

$$\text{Entropy_Info}(T) = \text{Entropy_Info}(3, 1) =$$

$$= -\left[\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right]$$

$$= -(-0.3111 + -0.4997)$$

$$= 0.8108$$

$$\text{Entropy_Info}(T, \text{Interactiveness}) = \frac{2}{4} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right] + \frac{2}{4} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$= 0 + 0.4997$$

$$\text{Gain}(\text{Interactiveness}) = 0.8108 - 0.4997$$

$$= 0.3111$$

$$\text{Entropy_Info}(T, \text{Practical Knowledge})$$

$$= \frac{2}{4} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right] + \frac{1}{4} \left[-\frac{0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1} \right] + \frac{1}{4} \left[-\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= 0$$

$$\text{Gain}(\text{Practical Knowledge}) = 0.8108$$

$$\text{Entropy_Info}(T, \text{Communication Skills})$$

$$= \frac{2}{4} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right] + \frac{1}{4} \left[-\frac{0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1} \right] + \frac{1}{4} \left[-\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} \right]$$

$$= 0$$

$$\text{Gain}(\text{Communication Skills}) = 0.8108$$

The gain calculated for all the attributes is shown in Table 6.9.

Table 6.9: Total Gain

Attributes	Gain
Interactiveness	0.3111
Practical Knowledge	0.8108
Communication Skills	0.8108

Here, both the attributes 'Practical Knowledge' and 'Communication Skills' have the same Gain. So, we can either construct the decision tree using 'Practical Knowledge' or 'Communication Skills'. The final decision tree is shown in Figure 6.4.

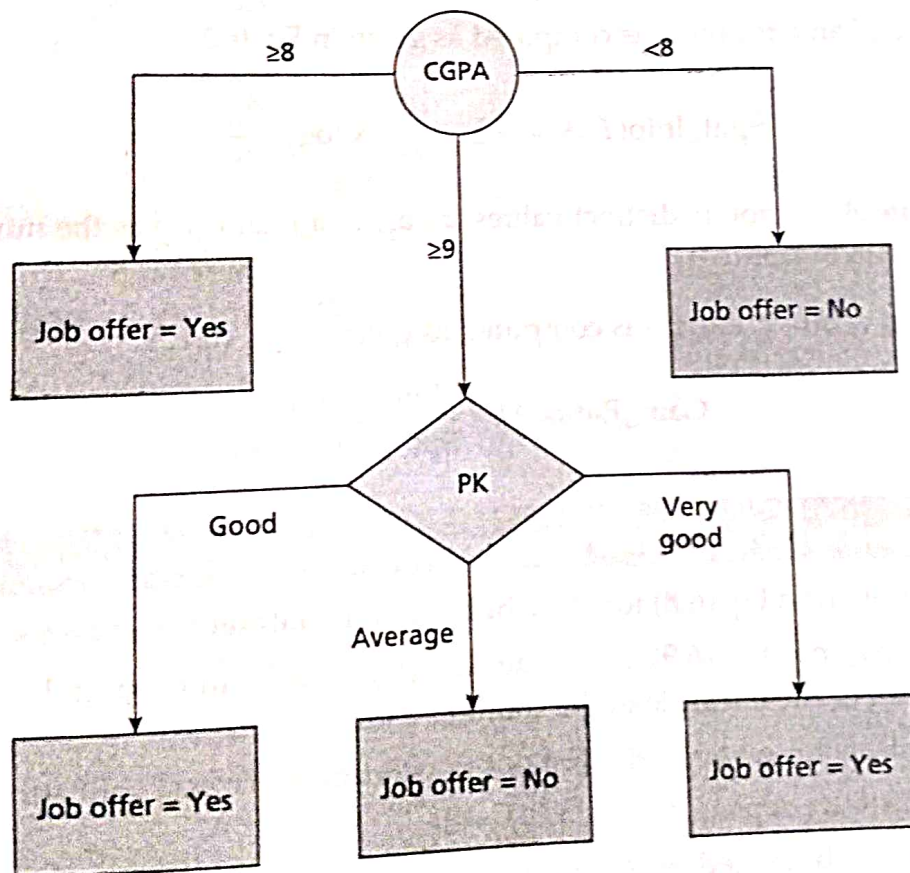


Figure 6.4: Final Decision Tree